# EXHIBIT 4

```
 1              UNITED STATES DISTRICT COURT
 2               DISTRICT OF MASSACHUSETTS
 3
 4
 5    SINGULAR COMPUTING LLC,          )
                                       )
 6           Plaintiff,                )
                                       )
 7       vs.                           ) Case Nos.
                                       ) 1:19-cv-12551-FDS
 8    GOOGLE LLC,                      )
                                       )
 9           Defendant.                )
                                       )
10
11
12
13
14    ***  ███████████████████████████  ***
15
16
17              REMOTE VIDEO DEPOSITION OF
18              DR. SUNIL KHATRI, VOLUME II
19
20
21
22
23    DATE TAKEN:  MARCH 24, 2023
24    REPORTED BY:  RENEE HARRIS, CSR 14168, CCR, RPR
      JOB NO. 5805112
25    PAGES:  350 - 699
```

Page 350

1        A.   Okay.  I'm at the claims section.

2        Q.   Yeah, and you understand that it's Claim

3    7 that's being -- of the '156 patent that's being

4    asserted in this case; right?

5        A.   Yes, I'd like to just double-check that        05:32:11

6    to be 100 percent sure.

7        Q.   Sure.

8        A.   Yes, it is Claim 7 of the '156 patent;

9    and that's the dependent claim which depends on 3,

10   which depends on 2, which depend on 1.                   05:32:33

11       Q.   Claim 1 has been ruled invalid; correct?

12           MR. SEEVE:  Objection.  Calls for

13       speculation.

14           THE WITNESS:  I have no idea about that.

15       I have no information to answer that question        05:32:50

16       either way.

17   BY MR. KAMBER:

18       Q.   You participated in the IPR proceedings

19   in this case; correct?

20       A.   I -- I did -- you know, I did present --        05:33:00

21   I mean, I did present -- I was deposed as well in

22   the IPR proceedings.  But the outcome of the IPR

23   proceedings, I'm unaware of.

24       Q.   You have no idea what the outcome of the

25   IPR proceedings is?                                      05:33:17

Page 555

1          That -- and that last part, the, "And

2     produces unexpected results," is an important

3     element that shouldn't be missed because --

4     because as Dr. Leeser says in her report that, oh,

5     there's nothing surprising about -- you know,          06:25:04

6     there's no surprising results you might get by

7     using additional execution units and so on.  But

8     this -- this line underscores that, no, there is a

9     surprising result that you obtain because it does

10    give unexpected results and these unexpected          06:25:21

11    results are what are described by, in the

12    specification of the patent, from -- you know,

13    from the columns I just recited to you which begin

14    around, you know, column 16, lines 59, all the way

15    through column 23 lines, I guess it was 34 or         06:25:39

16    something like this.  There's explicit disclosure

17    in the patent that -- of this unexpected result

18    and of course this unexpected result is also

19    described in Dr. Bates slides to Google as well.

20          So this line is actually underscoring          06:26:00

21    that unexpected result and it's important to

22    not -- to not leave out that fragment of the line

23    when you're citing it.

24       Q.  Dr. Khatri, is it your belief that

25    Dr. Bates was the first one to determine that        06:26:17

                                                    Page 578

1    using narrower bitwidths led to more parallelism?

2         MR. SEEVE:  Objection -- objection.

3       Mischaracterizes the report and the witness's

4       prior testimony.  Calls for a legal

5       conclusion.                                    06:26:31

6         THE WITNESS:  I'm going to say that, you

7       know, Dr. Bates -- you know, the question you

8       are asking is an incomplete question.

9       Because what I'm -- my answer is, Dr. Bates

10      was the first to describe multiple things and   06:26:45

11      among these things is, one, is that if you

12      use low precision, high-dynamic execution

13      units, A, you can fit more units, you know,

14      in the same circuit area; B, that results in

15      reduction -- I mean, B, that results in         06:27:00

16      dramatic improvement in performance at high

17      precision, that's important, right.  And

18      that's the surprising result.

19         That's the part that -- you know, the

20      totality of all these comments is what is       06:27:13

21      important, and that's the -- that's what is

22      described in the patent in detail, as well as

23      in Dr. Bates presentations to Google, which I

24      cite on page 53 of my report where he

25      shows -- you know, one thing he shows is        06:27:31

                                            Page 579

1       that, you know, the approximate

2       floating-point units are much smaller; and

3       therefore, the next citation is the figure at

4       the bottom of page 53, which comes from his

5       slides, which says that, you know, that we          06:27:44

6       can have, you know, 100x more floating-point

7       units compared to the IEEE floating-point

8       units.  That's the other thing he says.  That

9       means you can squeeze in more floating-point

10      units in the same chip area.                        06:27:59

11          But then the next slide which -- which I

12      cite in page 54 -- it shows that the software

13      can get, you know, 10,000x better speed and

14      power than the GPU, that's what the other

15      citation is.                                        06:28:15

16          And then finally, the other comment is

17      that, you know, he shows the -- you know,

18      which I -- which I show from his slide and

19      surprise No. 2, page 57 of my report, that

20      even though we have these low precision            06:28:28

21      units, you know, operating, you know, in

22      parallel, and I'll quote here this -- because

23      this is important.

24          It says, "The high precision CPU managing

25      low precision workers" -- that means LPHDR         06:28:42

Page 580

```
 1        execution units -- "can yield high precision

 2        results, like the CPU," completely unexpected

 3        result which is surprising and many Google

 4        engineers in their e-mail responses among

 5        each other expressed significant surprise at    06:28:57

 6        this, and also it says, surprise No. 2

 7        continues.

 8            It says, "but with size, power, cost of

 9        the low precision hardware for varied tasks."

10            So not only are we going to get this, you    06:29:10

11        know, size, power, and cost comparable to --

12        of the low precision hardware but also this

13        applies to many tasks.  This is significant

14        because this allows this idea to be used for

15        many tasks and get tremendous speed up, you    06:29:26

16        know -- you know, with these low precision

17        units, and the precision still is comparable

18        to the "high precision CPU," completely

19        surprising result.

20   BY MR. KAMBER:                                      06:29:39

21       Q.  Those ideas were known before, though;

22   right, Dr. Khatri?

23           MR. SEEVE:  Objection -- objection.

24       Mischaracterizes the witness's testimony.

25       Argumentative.                                  06:29:46
```

Page 581

1      THE WITNESS:  They were definitely not

2      known and there is basically -- that's the

3      inventiveness of the patent and that's also

4      expressed in the e-mails the engineers

5      exchanged among themselves once they saw the      06:30:02

6      second doc of Dr. Bates, and there was some

7      significant praise that they expressed,

8      significant surprise that they expressed.

9         There's multiple reasons why this was

10     surprising to the community, because the      06:30:15

11     conventional wisdom -- in fact, the patent

12     specification says this:  The conventional

13     wisdom is that if you want a high precision

14     algorithmic output, you must use high

15     precision execution units.      06:30:31

16         But this patent shows a completely

17     surprising result, that if you want -- if

18     you -- if you use low precision execution

19     units and you can use many of them, for many

20     applications, you can still get a high      06:30:44

21     precision output, which is significantly

22     surprising and it's completely against the

23     conventional wisdom in the field of -- in the

24     field.

25         And there is disclosure in the patent, I      06:30:57

                                        Page 582

```
 1          can point you to it, where -- where

 2          Dr. Bates, the inventor sort of describes

 3          this.

 4     BY MR. KAMBER:

 5          Q.  Let move to Exhibit 12 for a moment --    06:31:09

 6     Dr. Khatri, go to Exhibit 12, again, please.

 7              MR. SEEVE:  I'd like to point out that I

 8          think you just interrupted the witness's

 9          answer but --

10              MR. KAMBER:  I disagree.                  06:31:17

11              MR. SEEVE:  Like you've done so many

12          times.

13              MR. KAMBER:  He was done with his answer.

14          He was offering to --

15              THE WITNESS:  Let me open Exhibit 12 real 06:31:28

16          quick.  Excuse me.  I have Exhibit 12 open in

17          front of me.

18     BY MR. KAMBER:

19          Q.  Go to page 5, please.

20          A.  Can you give me, if you don't mind, the   06:31:44

21     title of that?

22          Q.  "Format design trade-offs."

23          A.  I see that slide.

24          Q.  This slide shows, as Dr. Leeser was

25     explaining, at this HPEC conference, that using    06:32:01
```

Page 583

1        which means you must use wider bitwidth,

2        which means you should get a high precision

3        functional unit, not a low precision

4        execution unit.  The surprising -- let me

5        please finish -- the surprising and            06:33:19

6        significant aspect of the asserted patents is

7        that despite using narrower bitwidths, you

8        can get a high precision output which

9        doctor -- you know, which the patent

10       describes -- which Dr. Bates describes in      06:33:36

11       those paragraphs that I cited to you which I

12       think were paragraphs -- sorry, columns 14

13       through columns 23.

14           Those are concrete examples and concrete

15       experiments that Dr. Bates had conducted to    06:33:49

16       show that with narrower bitwidths, one can

17       still get high precision for many

18       applications, and that's completely

19       contradictory to this slide because this

20       slide says, to get high precision, you should  06:34:03

21       use wider bitwidths, because as the arrow

22       pointing to the right, saying wider bitwidths

23       gives rise to higher precision.

24           So Dr. Bates' observation, Dr. Bates'

25       patent and the asserted claims and the         06:34:17

Page 585

1      asserted patents show this completely

2      surprising phenomena, which the conventional

3      wisdom, you know, simply didn't subscribe to,

4      which is why, as I said, you know, ████████

█      ████████████████████████████████        ███████████

█      ████████████████████████████

█          ██████████████████████████████

█      ██████████████████████████████████

█      ████████████████████████████

██      ██████████████████████████        ███████████

██      ██████████████████████████

██      ██████████████████████████

██      ██████████████████████████

██      ████████████████████████████████

██      ████████████████████████        ███████████

16     BY MR. KAMBER:

17         Q.   In that response, are you referring to

18     low precision as construed by the Court or in some

19     other sense?

20         MR. SEEVE:   Objection.  Mischaracterizes          06:35:11

21     the witness's testimony.  Vague and

22     ambiguous.

23         THE WITNESS:   I don't understand your

24     question, so -- when you say, when I'm

25     referring to "precision," what do you mean?          06:35:22

Page 586

1           that -- sorry.

2               For all the IPRs, because we are talking

3           about IPRs, for all the IPRs that were --

4           that were filed, I'm not aware of the -- the

5           legal paperwork that's filed back-and-forth        07:21:35

6           between Google and -- you know, and the PTAB

7           or -- or Singular and the PTAB or Google and

8           Singular.

9               I'm only aware of those documents that

10          were made available to me for the analysis         07:21:51

11          that I needed to conduct which is purely

12          technical and of course not legal because I'm

13          not a lawyer.

14              So whatever documents were provided to me

15          for my technical analyses, which were all          07:22:01

16          that I requested, those I reviewed.  But

17          subsequent to that, I'm not -- I'm not aware

18          of the rulings or decisions that the PTAB has

19          made about specific claim elements.

20              And to the extent that their analysis          07:22:18

21          defers from mine, I respectfully accept it

22          but I disagree with it because I stand by my

23          analysis.

24     BY MR. KAMBER:

25          Q.  Do you have any understanding that             07:22:29

                                                    Page 612